

Fountain Codes Based Distributed Storage Algorithms for Large-scale Wireless Sensor Networks

Salah A. Aly
Dept. of Computer Science
Texas A&M University
College Station, TX 77843
salah@cs.tamu.edu

Zhenning Kong
Dept. of Electrical Engineering
Yale University
New Haven, CT 06520
zhenning.kong@yale.edu

Emina Soljanin
Bell Laboratories
Alcatel-Lucent
Murray Hill, NJ 07974
emina@lucent.com

Abstract

We consider large-scale networks with n nodes, out of which k are in possession, (e.g., have sensed or collected in some other way) k information packets. In the scenarios in which network nodes are vulnerable because of, for example, limited energy or a hostile environment, it is desirable to disseminate the acquired information throughout the network so that each of the n nodes stores one (possibly coded) packet and the original k source packets can be recovered later in a computationally simple way from any $(1 + \epsilon)k$ nodes for some small $\epsilon > 0$.

We developed two distributed algorithms for solving this problem based on simple random walks and Fountain codes. Unlike all previously developed schemes, our solution is truly distributed, that is, nodes do not know n , k or connectivity in the network, except in their own neighborhoods, and they do not maintain any routing tables. In the first algorithm, all the sensors have the knowledge of n and k . In the second algorithm, each sensor estimates these parameters through the random walk dissemination. We present analysis of the communication/transmission and encoding/decoding complexity of these two algorithms, and provide extensive simulation results as well¹.

1 Introduction

Wireless sensor networks consist of small devices (sensors) with limited resources (e.g., low CPU power, small bandwidth, limited battery and memory). They can be deployed to monitor objects, measure temperature, detect fires, and other disaster phenomena. They are often used in

isolated, hard to reach areas, where human involvement is limited. Consequently, data acquired by sensors may have short lifetime, and any processing on it within the network should have low complexity and power consumption [18].

We consider a large-scale wireless sensor networks with n sensors. Among them, $k \ll n$ sensors have collected (sensed) some information. Since sensors are often short-lived because of limited energy or hostile environment, it is desirable to disseminate the acquired information throughout the network so that each of the n nodes stores one (possibly coded) packet and the original k source packets can be recovered in a computationally simple way from any $(1 + \epsilon)k$ of nodes for some small $\epsilon > 0$. Here, the sensors do not know locations of each other, and they do not maintain any routing tables.

Various solutions to the centralized version of this problem have been proposed, and are based on well known coding schemes such as Fountain codes [6] or MDS codes [16]. To distribute the information from multiple sources throughout the network so that each node stores a coded packet as if obtained by centralized LT (Luby Transform) coding [12], Lin *et al.* [11] proposed a solution that uses random walks with traps. To achieve the desired code degree distribution, they employed the Metropolis algorithm to specify transition probabilities of the random walks. In this way, the original k source packets are encoded by LT codes and the decoding process can be done by querying any $(1 + \epsilon)k$ arbitrary sensors. Because of properties of LT codes, the encoding and decoding complexity are linear and therefore have low energy consumption.

In the methods of [11], the knowledge of the total number of sensors n and sources k is required for calculating the number of random walks that each source needs to initiate and for calculating the probability of trapping at each sensor. Another type of global information, namely, the maximum node degree (i.e., the maximum number of neighbors) in the network, is also required to perform the Metropolis

¹This work was accomplished while S.A.A and Z.K. were spending a summer research internship at Bell Labs & Alcatel-Lucent, Murray Hill, N.J., 2007, and it was submitted as US patent in [2]. They would like to thank Bell Labs & Alcatel-Lucent staff members for their hospitality.

algorithm. However, for a large-scale sensor network, such global information may not be easy to obtain by each individual sensor, especially when there is possibility of change in topology. Moreover, the algorithms proposed in [11] assume that each sensor encodes only after receiving enough source packets. This requires each sensor to maintain a large enough temporary memory buffer, which may not be practical in real sensor networks.

In this paper, we propose two new algorithms to solve the distributed storage problem in large-scale sensor networks. We refer to these algorithms as LT-Codes based Distributed Storage-I (LTCDS-I) and LT-Codes based Distributed Storage-II (LTCDS-II). Both algorithms use simple random walks without trapping to disseminate source packets. In contrast to the methods in [11], both algorithms demand little global information and memory at each sensor. In LTCDS-I, only the values of n and k are needed, whereas the maximum node degree, which is more difficult to obtain, is not required. In LTCDS-II, no sensor needs to know any global information (that is, knowing n and k is no longer required). Instead, sensors can obtain good estimates for those parameters by using some properties of random walks. Moreover, in both algorithms, instead of waiting until all the necessary source packets are collected to do encoding, each sensor makes decisions and performs encoding online upon each reception of resource packets. This mechanism reduces the memory demand significantly.

The main contributions of this paper are as follows:

- (i) We propose two new algorithms (LTCDS-I and LTCDS-II) for distributed storage in large-scale sensor networks, using simple random walks and LT codes. These algorithms are simpler, more robust, and less constrained in comparison to previous solutions.
- (ii) We present complexity analysis of both algorithms, including transmission, encoding, and decoding complexity.
- (iii) We evaluate and illustrate the performance of both algorithms by extensive simulation.

This paper is organized as follows. We start with a short survey of the related work in Section 2. In Section 3, we introduce the network model and present Luby Transform (LT) codes. In Section 4, we propose two LT codes based distributed storage algorithms called LTCDS-I and LTCDS-II. We then present simulation studies and provide performance analysis of the proposed algorithms in Section 5, and concluded in Section 6.

2 Related Work

The most related work to one presented here is [11, 10]. Lin *et al.* studied the question “how to retrieve historical data that the sensors have gathered even if some sensors are destroyed or disappeared from the network?” They ana-

lyzed techniques to increase *persistence* of sensed data in a random wireless sensor network, and proposed two decentralized algorithms using Fountain codes to guarantee the persistence and reliability of cached data on unreliable sensors. They used random walks to disseminate data from multiple sensors (sources) to the whole network. Based on the knowledge of the total number of sensors n and sources k , each source calculates the number of random walks it needs to initiate, and each sensor calculates the number of source packets it needs to trap. In order to achieve some desired packet distribution, the transition probabilities of random walks are specified by the well known Metropolis algorithm [11].

Dimakis *et al.* in [4, 6] proposed a decentralized implementation of Fountain codes that uses geographic routing, where every node has to know its location. The motivation for using Fountain codes is their low decoding complexity. Also, one does not know in advance the degrees of the output nodes in this type of codes. The authors proposed a randomized algorithm that constructs Fountain codes over a grid network using only geographical knowledge of nodes and local randomized decisions. Fast random walks are used to disseminate source data to the storage nodes in the network.

Kamara *et al.* in [9, 8] proposed a novel technique called *growth codes* to increase data *persistence* in wireless sensor networks, namely, increase the amount of information that can be recovered at the sink. Growth coding is a linear technique in which information is encoded in an online distributed way with increasing degree of a storage node. Kamara *et al.* showed that *growth codes* can increase the amount of information that can be recovered at any storage node at any time period whenever there is a failure in some other nodes. They did not use robust or soliton distributions, but proposed a new distribution depending on the network condition to determine degrees of the storage nodes. The motivation for their work was that i) Positions and topology of the nodes are not known. ii) They assume a round time of node updates, meaning with increasing the time t , degree of a symbol is increased. This is the idea behind growth degrees. iii) They provide practical implementations of growth codes and compare its performance with other codes. iv) The decoding part is done by querying an arbitrary sink, if the original sensed data has been collected correctly then finish, otherwise query another sink node.

Lun *et al.* in [13] proposed two decentralized algorithms to compute the minimum-cost subgraphs for establishing multicast connections using network coding. Also, they extended their work to the problem of minimum-energy multicast in wireless networks as well as they studied directed point-to-point multicast and evaluated the case of elastic rate demand.

3 Wireless Sensor Networks and Fountain Codes

In this section, we introduce our network model and provide background of Fountain codes and, in particular, one important class of Fountain codes—LT (Luby Transform) codes [12].

3.1 Network Model

Our wireless sensor network consists of n nodes that are uniformly distributed at random in a region $\mathcal{A} = [L, L]^2$ for $L > 1$. The *density* of the network is given by

$$\lambda = \frac{n}{|\mathcal{A}|} = \frac{n}{L^2}, \quad (1)$$

where $|\mathcal{A}|$ is the two-dimensional Lebesgue measure (or area) of \mathcal{A} . Each sensor node has an identical communication radius 1; thus any two nodes can communicate with each other if and only if their distance is less than or equal to 1. This model is known as random geometric graphs [7, 15]. Among these n nodes, there are k source nodes that have information to be disseminated throughout the network for storage. These k nodes are uniformly and independently distributed at random among the n nodes. Usually, the fraction of source nodes, i.e., $\frac{k}{n}$, is not very large (e.g., 10%, or 20%).

Note that, although we assume the nodes are uniformly distributed at random in a region, our algorithms and results do not rely on this assumption. In fact, they can be applied for any network topology, for example, regular grids.

We assume that no node has knowledge about the locations of other nodes and no routing table is maintained; consequently, the algorithm proposed in [5] cannot be applied. Moreover, we assume that each node has limited or no knowledge of global information, but know its neighbors. The limited global information refers to the total numbers of nodes n and sources k . Any further global information, for example the maximal number of neighbors in the network, is not available. Hence, the algorithms proposed in [11, 10] are not applicable.

Definition 1. (Node Degree) Consider a graph $G = (V, E)$, where V and E denote the set of nodes and links, respectively. Given $u, v \in V$, we say u and v are adjacent (or u is adjacent to v , and vice versa) if there exists a link between u and v , i.e., $(u, v) \in E$. In this case, we also say that u and v are neighbors. Denote by $\mathcal{N}(u)$ the set of neighbors of a node u . The number of neighbors of a node u is called the node degree of u , and denoted by $d_n(u)$, i.e., $|\mathcal{N}(u)| = d_n(u)$. The mean degree of a graph G is then given by

$$\mu = \frac{1}{|V|} \sum_{u \in G} d_n(u), \quad (2)$$

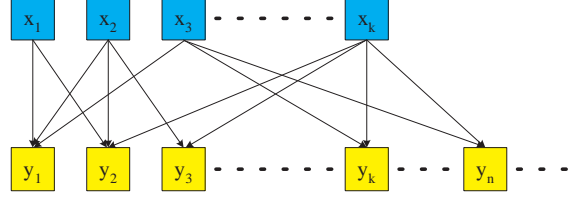


Figure 1. The encoding operations of Fountain codes: each output is obtained by XOR-ing d source blocks chosen uniformly and independently at random from k source inputs, where d is drawn according to a probability distribution $\Omega(d)$.

where $|V|$ is the total number of nodes in G .

3.2 Fountain Codes

For k source blocks $\{x_1, x_2, \dots, x_k\}$ and a probability distribution $\Omega(d)$ with $1 \leq d \leq k$, a Fountain code with parameters (k, Ω) is a potentially limitless stream of output blocks $\{y_1, y_2, \dots\}$. Each output block is obtained by XORing d randomly and independently chosen source blocks, where d is drawn from a specially designed distribution $\Omega(d)$. This is illustrated in Figure 3.2. Fountain codes are rateless, and one of their main advantage is that the encoding operations can be performed online. The encoding cost is the expected number of operation sufficient for generating an output symbol, and the decoding cost is the expected number of operations sufficient to recover the k input blocks. Another advantage of Fountain codes, as opposed to purely random codes is that their decoding complexity can be made low by appropriate choice of $\Omega(d)$, with little sacrifice in performance. The decoding of Fountain codes can be done by message passing.

Definition 2. (Code Degree) For Fountain codes, the number of source blocks used to generate an encoded output y is called the code degree of y , and denoted by $d_c(y)$. By construction, the code degree distribution $\Omega(d)$ is the probability distribution of $d_c(y)$.

3.3 LT Codes

LT (Luby Transform) codes are a special class of Fountain codes which uses *Ideal Soliton* or *Robust Soliton* distributions [12]. The Ideal Soliton distribution $\Omega_{is}(d)$ for k source blocks is given by

$$\Omega_{is}(i) = \Pr(d = i) = \begin{cases} \frac{1}{k}, & i = 1 \\ \frac{1}{i(i-1)}, & i = 2, 3, \dots, k. \end{cases} \quad (3)$$

Let $R = c_0 \sqrt{k} \ln(k/\delta)$, where c_0 is a suitable constant and $0 < \delta < 1$. The Robust Soliton distribution for k source blocks is defined as follows. Define

$$\tau(i) = \begin{cases} \frac{R}{ik}, & i = 1, \dots, \frac{k}{R} - 1 \\ \frac{R \ln(R/\delta)}{k}, & i = \frac{k}{R}, \\ 0, & i = \frac{k}{R} + 1, \dots, k, \end{cases} \quad (4)$$

and let

$$\beta = \sum_{i=1}^k \tau(i) + \Omega_{is}(i). \quad (5)$$

The Robust Soliton distribution is given by

$$\Omega_{rs}(i) = \frac{\tau(i) + \Omega_{is}(i)}{\beta}, \text{ for all } i = 1, 2, \dots, k \quad (6)$$

The following result provides the performance of the LT codes with Robust Soliton distribution [12, Theorems 12 and 13].

Lemma 3 (Luby [12]). *For LT codes with Robust Soliton distribution, k original source blocks can be recovered from any $k + O(\sqrt{k} \ln^2(k/\delta))$ encoded output blocks with probability $1 - \delta$. Both encoding and decoding complexity is $O(k \ln(k/\delta))$.*

4 LT-Codes Based Distributed Storage (LTCDS) Algorithms

In this section, we present two LT-Codes based Distributed Storage (LTCDS) algorithms. In both algorithms, the source packets are disseminated throughout the network by a simple random walk. In the first one, called LTCDS-I algorithm, we assume that each node in the network has limited the global information, that is, knows the total number of sources k and the total number of nodes n . Unlike the scheme proposed in [10], our algorithm does not require the nodes to know the maximum degree of the graph, which is much harder to obtain than k and n . The second algorithm, called LTCDS-II, is a fully distributed algorithm which does not require nodes to know any global information. The price we pay for this benefit is extra transmissions of the source packets to obtain estimates for n and k .

4.1 With Limited Global Information—LTCDS-I

In LTCDS-I, we assume that each node in the network knows the values of k and n . We use simple random walks [1, 17] for each source to disseminate its information to the whole network. At each round, each node u that has

packets to transmit chooses one node v among its neighbors uniformly independently at random, and sends the packet to the node v . In order to avoid local-cluster effect—each source packet is trapped most likely by its neighbor nodes—we let each node accept a source packet equiprobably. To achieve this, we also need each source packet to visit each node in the network at least once.

For a random walk on a graph, the *cover time* is defined as follows [1, 17]:

Definition 4. (Cover Time) *Given a graph G , let $T_{cover}(u)$ be the expected length of a random walk that starts at node u and visits every node in G at least once. The cover time of G is defined by*

$$T_{cover}(G) = \max_{u \in G} T_{cover}(u). \quad (7)$$

For a simple random walk on a random geometric graph, the following result bounds the cover time [3].

Lemma 5 (Avin and Ercal [3]). *If a random geometric graph with n nodes is a connected graph with high probability, then*

$$T_{cover}(G) = \Theta(n \log n). \quad (8)$$

As a result of Lemma 5, we can set a counter for each source packet and increase the counter by one after each forward transmission until the counter reaches some threshold $C_1 n \log n$ to guarantee that the source packet visits each node in the network at least once. The detailed descriptions of the initialization, encoding and storage phases (steps) of LTCDS-I algorithm are given below:

(i) Initialization Phase:

- (1) Each node u in the network draws a random number $d_c(u)$ according to the distribution $\Omega_{is}(d)$ given by (3) (or $\Omega_{rs}(d)$ given by (6)). Each source node $s_i, i = 1, \dots, k$ generates a header for its source packet x_{s_i} and puts its ID and a counter $c(x_{s_i})$ with initial value zero into the packet header. We set up tokens for initial and update packets. We assume that a token is set to zero for an initial packet and 1 for an update packet.

$$packet_{s_i} = (ID_{s_i}, x_{s_i}, c(x_{s_i}))$$

- (2) Each source node s_i sends out its own source packet x_{s_i} to another node u which is chosen uniformly at random among all its neighbors $\mathcal{N}(s_i)$.
- (3) The chosen node u accepts this source $packet_{s_i}$ with probability $\frac{d_c(u)}{k}$ and updates its storage as

$$y_u^+ = y_u^- \oplus x_{s_i}, \quad (9)$$

where y_u^- and y_u^+ denote the packet that the node u stores before and after the updating, respectively, and \oplus represents XOR operation. No matter whether the source packet is accepted or not,

the node u puts it into its forward queue and set the counter of x_{s_i} as

$$c(x_{s_i}) = 1. \quad (10)$$

(ii) **Encoding Phase:**

- (1) In each round, when a node u receives at least one source packet before the current round, u forwards the head-of-line (HOL) packet x in its forward queue to one of its neighbor v , chosen uniformly at random among all its neighbors $\mathcal{N}(u)$.
- (2) Depending on how many times x has visited v , the node v makes its decisions:
 - If it is the first time that x visits v , then the node v accepts this source packet with probability $\frac{d}{k}$ and updates its storage as

$$y_v^+ = y_v^- \oplus x. \quad (11)$$

- If x has visited v before and $c(x) < C_1 n \log n$ where C_1 is a system parameter, then the node v accepts this source packet with probability 0.
- No matter x is accepted or not, the node v puts it into its forward queue and increases the counter of x by one:

$$c(x) = c(x) + 1. \quad (12)$$

- If x has visited v before and $c(x) \geq C_1 n \log n$ then the node v discards the packet x forever.

(iii) **Storage Phase:**

When a node u makes its decisions for all the source packets $x_{s_1}, x_{s_2}, \dots, x_{s_k}$, i.e., all these packets have visited the node u at least once, the node u finishes its encoding process by declaring the current y_u to be its storage packet.

The pseudo-code of these steps is given in LTCDS-I Algorithm 1.

The following theorem establishes the code degree distribution of each storage node induced by the LTCDS-I algorithm.

Theorem 6. *When a sensor network with n nodes and k sources finishes the storage phase of the LTCDS-I algorithm, the code degree distribution of each storage node u is given by*

$$\begin{aligned} & \Pr(\tilde{d}_c(u) = i) \\ &= \sum_{d_c(u)=1}^k \binom{k}{i} \left(\frac{d_c(u)}{k} \right)^i \left(1 - \frac{d_c(u)}{k} \right)^{k-i} \Omega'(d_c(u)), \end{aligned} \quad (13)$$

where $d_c(u)$ is given in the initialization phase of the LTCDS-I algorithm from distribution $\Omega'(d)$ (i.e., $\Omega_{is}(d)$ or $\Omega_{rs}(d)$), and $\tilde{d}_c(u)$ is the code degree of the node u resulting from the algorithm.

Input: number of nodes n , number of sources k , source packets $x_{s_i}, i = 1, 2, \dots, k$ and a positive constant C_1

Output: storage packets $y_i, i = 1, 2, \dots, n$

foreach node $u = 1 : n$ **do**

 Generate $d_c(u)$ according to $\Omega_{is}(d)$ (or $\Omega_{rs}(d)$);

end

foreach source node $s_i, i = 1 : k$ **do**

 Generate header of x_{s_i} and $token = 0$;

$c(x_{s_i}) = 0$;

 Choose $u \in \mathcal{N}(s_i)$ uniformly at random, send x_{s_i} to u ;

$coin = \text{rand}(1)$;

if $coin \leq \frac{d_c(u)}{k}$ **then** $y_u = y_u \oplus x_{s_i}$;

 Put x_{s_i} into u 's forward queue;

$c(x_{s_i}) = c(x_{s_i}) + 1$;

end

while source packets remaining **do**

foreach node u receives packets before current round **do**

 Choose $v \in \mathcal{N}(u)$ uniformly at random;

 Send HOL packet x_{s_i} in u 's forward queue to v ;

if v receives x_{s_i} for the first time **then**

$coin = \text{rand}(1)$;

if $coin \leq \frac{d_c(v)}{k}$ **then**

$y_v = y_v \oplus x_{s_i}$;

 Put x_{s_i} into v 's forward queue;

$c(x_{s_i}) = c(x_{s_i}) + 1$

end

else if $c(x_{s_i}) < C_1 n \log n$ **then**

 Put x_{s_i} into v 's forward queue;

$c(x_{s_i}) = c(x_{s_i}) + 1$;

else

 Discard x_{s_i} ;

end

end

end

Algorithm 1: LTCDS-I Algorithm: LT-Codes based Distributed Storage Algorithm for a wireless sensor network (WSN) with limited global information, i.e., values of n and k are known at every node. It consists of three phases: initialization, encoding and storage phases. The algorithm can also be deployed in a WSN after estimating values of n and k , as shown in LTCDS-II algorithm.

Proof. For each node u , $d_c(u)$ is drawn from a distribution $\Omega'(d)$ (i.e., $\Omega_{is}(d)$ or $\Omega_{rs}(d)$). Given $d_c(u)$, the node u accepts each source packet with probability $\frac{d_c(u)}{k}$ independently of each other and $d_c(u)$. Thus, the number of source packets that the node u accepts follows a Binomial distribu-

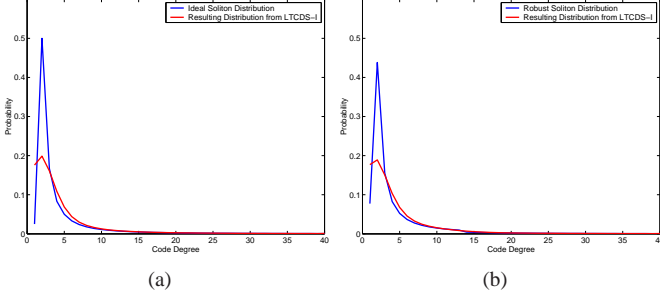


Figure 2. Code degree distribution comparing: (a) Ideal Soliton distribution Ω_{is} (given by (3)) and the resulting degree distribution from LTCDS-I algorithm (given by (13)). Here $k = 40$; (b) Robust Soliton distribution Ω_{rs} (given by (6)) and the resulting degree distribution from LTCDS-I algorithm (given by (13)). Here $k = 40$, $c_0 = 0.1$ and $\delta = 0.5$.

tion with parameter $\frac{d_c(u)}{k}$. Hence,

$$\begin{aligned} \Pr(\tilde{d}_c(u) = i) &= \sum_{d_c(u)=1}^k \Pr(\tilde{d}_c(u) = i | d_c(u)) \Omega'(d_c(u)) \\ &= \sum_{d_c(u)=1}^k \binom{k}{i} \left(\frac{d_c(u)}{k} \right)^i \left(1 - \frac{d_c(u)}{k} \right)^{k-i} \Omega'(d_c(u)), \end{aligned}$$

and thereafter (13) holds. \square

Theorem 6 indicates that the code degree $\tilde{d}_c(u)$ is not the same as $d_c(u)$. In fact, one may achieve the exact desired code degree distribution by letting all the sensors hold the received source packets in their temporary buffer until they collect all k source packets. Then they can randomly choose $d_c(u)$ packets. In this way, the resulting degree distribution is exactly the same as Ω_{is} or Ω_{rs} . However, this requires that each sensor has enough buffer or memory, which is usually not practical, especially when k is large. Therefore, in LTCDS-I, we assume each sensor has very limited memory and let them make their decision upon each reception.

Fortunately, from Figure 2, we can see that at the high degree end, the resulting code degree distribution obtained by the LTCDS-I algorithm (13) perfectly matches the desired code degree distribution, i.e., either the Ideal Soliton distribution Ω_{is} (3) or the Robust Soliton distribution Ω_{rs} (6). For the resulting degree distribution and the desired degree distributions, the difference only lies at the low degree end, especially at degree 1 and degree 2. In particular, the resulting degree distribution has higher probability at degree 1 and lower probability at degree 2 than the desired de-

gree distributions. The fact that higher probability at degree 1 turns out to compensate the lower probability at degree 2 so that the resulting degree distribution has very similar encoding and decoding behavior as LT codes using either the Ideal Soliton distribution or the Robust Soliton distribution. In our future study, we will provide theoretical analysis and prove that the degree distribution in 13 is equivalent, but not the same, as the degree distributed used in LT encoding [12]. Therefore, we have the following theorem, which can be proved by the same method for Lemma 3, see [12].

Theorem 7. Suppose sensor networks have n nodes and k sources and the LTCDS-I algorithm uses the Robust Soliton distribution Ω_{rs} . Then, when n and k are sufficient large, the k original source packets can be recovered from any $k + O(\sqrt{k} \ln^2(k/\delta))$ storage nodes with probability $1 - \delta$. The decoding complexity is $O(k \ln(k/\delta))$.

Theorem 7 asserts that when n and k are sufficiently large, the performance of the LTCDS-I is similar to LT coding.

Another main performance metric is the transmission cost of the algorithm, which is characterized by the total number of transmissions (the total number of steps of k random walks).

Theorem 8. Denote by $T_{LTCDS}^{(I)}$ the total number of transmissions of the LTCDS-I algorithm, then we have

$$T_{LTCDS}^{(I)} = \Theta(kn \log n), \quad (14)$$

where k is the total number of sources, and n is the total number of nodes in the network.

Proof. We know that each one of k source packets is stooped and discarded if and only if it has been forwarded for $C_1 n \log(n)$ times, for some constant C_1 . Then the total number of transmissions of the LTCDS-I algorithm for all k packets is a direct consequence and it is given by (14). \square

4.2 Without any Global Information—LTCDS-II

In many scenarios, especially when a change in network topology occurs because of, for example, node mobility or node failures, the exact values of n and k may not be available to all nodes. Therefore, to design a fully distributed storage algorithm which does not require any global information is very important and useful. In this subsection, we present such an algorithm based on LT codes, called LTCDS-II. The idea behind this algorithm is to utilize some features of simple random walks to do inference to obtain individual estimates of n and k for each node.

We introduce of *inter-visit time* and *inter-packet time* [1, 17, 14] as follows:

Definition 9. (*Inter-Visit Time*) For a random walk on a graph, the inter-visit time of node u , $T_{visit}(u)$, is the amount of time between any two consecutive visits of the random walk to node u . This inter-visit time is also called return time.

For a simple random walk on random geometric graphs, the following lemma provides results on the expected inter-visit time of any node. The proof is straightforward by following the standard result of stationary distribution of a simple random walk on graphs and the mean return time for a Markov chain [1, 17, 14]. For completeness, we provide the proof in Appendix 6.1.

Lemma 10. For a node u with node degree $d_n(u)$ in a random geometric graph, the mean inter-visit time is given by

$$E[T_{visit}(u)] = \frac{\mu n}{d_n(u)}, \quad (15)$$

where μ is the mean degree of the graph given by Equation (2).

From Lemma 10, we can see that if each node u can measure the expected inter-visit time $E[T_{visit}(u)]$, then the total number of nodes n can be estimated by

$$n = \frac{d_n(u)E[T_{visit}(u)]}{\mu}. \quad (16)$$

However, the mean degree μ is a global information and may be hard to obtain. Thus, we make a further approximation and let the estimate of n by the node u be

$$\hat{n}(u) = E[T_{visit}(u)]. \quad (17)$$

Hence, every node u computes its own estimate of n . In our distributed storage algorithms, each source packet follows a simple random walk. Since there are k sources, we have k individual simple random walks in the network. For a particular random walk, the behavior of the return time is characterized by Lemma 10. On the other hand, Lemma 12 below provides results on the inter-visit time among all k random walks, which is called inter-packet time for our algorithm, defined as follows:

Definition 11. (*Inter-Packet Time*) For k random walks on a graph, the inter-packet time of node u , $T_{packet}(u)$, is the amount of time between any two consecutive visits of those k random walks to node u .

For the mean value of inter-packet time, we have the following lemma, for which the proof is given in Appendix 6.2.

Lemma 12. For a node u with node degree $d_n(u)$ in a random geometric graph with k simple random walks, the mean inter-packet time is given by

$$E[T_{packet}(u)] = \frac{E[T_{visit}(u)]}{k} = \frac{\mu n}{k d_n(u)}, \quad (18)$$

where μ is the mean degree of the graph given by (2).

From Lemma 10 and Lemma 12, it is easy to see that for any node u , an estimation of k can be obtained by

$$\hat{k}(u) = \frac{E[T_{visit}(u)]}{E[T_{packet}(u)]}. \quad (19)$$

After obtaining estimates for both n and k , we can employ similar techniques used in LTCDS-I to do LT coding and storage. The detailed descriptions of the initialization, inference, encoding, and storage phases of LTCDS-II algorithm are given below:

(i) **Initialization Phase:**

- (1) Each source node $s_i, i = 1, \dots, k$ generates a header for its source packet x_{s_i} and puts its ID and a counter $c(x_{s_i})$ with initial value zero into the packet header.
- (2) Each source node s_i sends out its own source packet x_{s_i} to one of its neighbors u , chosen uniformly at random among all its neighbors $\mathcal{N}(s_i)$.
- (3) The node u puts x_{s_i} into its forward queue and sets the counter of x_{s_i} as

$$c(x_{s_i}) = 1. \quad (20)$$

(ii) **Inference Phase:**

- (1) For each node u , suppose $x_{s(u)_1}$ is the first source packet that visits u , and denote by $t_{s(u)_1}^{(j)}$ the time when $x_{s(u)_1}$ has its j -th visit to the node u . Meanwhile, each node u also maintains a record of visiting time for each other source packet $x_{s(u)_i}$ that visited it. Let $t_{s(u)_i}^{(j)}$ be the time when source packet $x_{s(u)_i}$ has its j -th visit to the node u . After $x_{s(u)_1}$ visiting the node u C_2 times, where C_2 is system parameter which is a positive constant, the node u stops this monitoring and recoding procedure. Denote by $k(u)$ the number of source packets that have visited at least once upon that time.
- (2) For each node u , let $J(s(u)_i)$ be the number of visits of source packet $x_{s(u)_i}$ to the node u and let

$$T_{s(u)_i} = \frac{1}{J(s(u)_i)} \sum_{j=1}^{J(s(u)_i)} t_{s(u)_i}^{(j+1)} - t_{s(u)_i}^{(j)} \quad (21)$$

$$= \frac{1}{J(s(u)_i)} (t_{s(u)_i}^{(J(s(u)_i))} - t_{s(u)_i}^{(1)}). \quad (22)$$

Then, the average inter-visit time for node u is given by

$$\bar{T}_{visit}(u) = \frac{1}{k(u)} \sum_{i=1}^{k(u)} T_{s(u)_i}. \quad (23)$$

Let $J_{min} = \min_{s(u)_i} \{t_{s(u)_i}^{(1)}\}$ and $J_{max} = \max_{s(u)_i} \{t_{s(u)_i}^{(J(s(u)_i))}\}$, then the inter-packet time

is given by

$$\bar{T}_{packet}(u) = \frac{J_{min} - J_{max}}{\sum_{s(u)_i} J(s(u)_i)}. \quad (24)$$

Then the node u can estimate the total number of nodes in the network and the total number of sources as

$$\hat{n}(u) = \bar{T}_{visit}(u), \quad (25)$$

and

$$\hat{k}(u) = \frac{\bar{T}_{visit}(u)}{\bar{T}_{packet}(u)}. \quad (26)$$

- (3) In this phase, the counter $c(x_{s_i})$ of each source packet $c(x_{s_i})$ is incremented by one after each transmission.

(iii) **Encoding Phase:**

When a node u obtains estimates $\hat{n}(u)$ and $\hat{k}(u)$, it begins encoding phase which is the same as the one in LTCDS-I Algorithm except that the code degree $d_c(u)$ is drawn from distribution $\Omega_{is}(d)$ (or $\Omega_{rs}(d)$) with replacement of k by $\hat{k}(u)$, and a source packet x_{s_i} is discarded if $c(x_{s_i}) \geq C_3 \hat{n}(u) \log \hat{n}(u)$, where C_3 is a system parameter which is a positive constant.

(iv) **Storage Phase:**

When a node u has made its decisions for \hat{k} source packets, it finishes its encoding process and y_u becomes the storage packet of u .

The total number of transmissions (the total number of steps of k random walks) in the LTCDS-II algorithm has the same order as LTCDS-I.

Theorem 13. Denote by $T_{LTCDS}^{(II)}$ the total number of transmissions of the LTCDS-II algorithm, then we have

$$T_{LTCDS}^{(II)} = \Theta(kn \log n), \quad (27)$$

where k is the total number of sources, and n is the total number of nodes in the network.

Proof. In the interference phase of the LTCDS-II algorithm, the total number of transmissions is upper bounded $C'n$ for some constants $C' > 0$. That is because each node needs to receive the first visit source packet for C_2 times, and by Lemma 10, the mean inter-visit time is $\Theta(n)$.

In the decoding phase, the same as in the LTCDS-I algorithm, in order to guarantee that each source packet visits all the nodes at least once, the number of steps of the simple random walk is $\Theta(n \log n)$. In other words, each source packet is stopped and discarded if and only if the counter reaches the threshold $C_3 n \log(n)$ for some system parameter C_3 . Therefore, we have (27). \square

4.3 Updating Data

Now, we turn our attention to data updating after all storage nodes saved their values y_1, y_2, \dots, y_n , but a sensor node, say s_i , wants to update its value to the appropriate set of storage nodes in the network. The following updating algorithm applies for both LTCDS-I and LTCDS-II. For simplicity, we illustrate the idea with LTCDS-I.

Assume the sensor node prepared a packet with its ID, old data x_{s_i} , new data x'_{s_i} along with a time-to-live parameter $c(s_i)$ initialized to zero. We will use also a simple random walk for data update.

$$packet_{s_i} = (ID_{s_i}, x_{s_i} \oplus x'_{s_i}, c(s_i)). \quad (28)$$

If we assume that the storage nodes keep ID's of the accepted packets, then the problem becomes simple. We just run a random walk and check for the coming packet's ID. Assume the node u keeps track of all ID's of its accepted packets. Then u accepts the updated message if ID of the coming packet is already included in the u 's ID list. Otherwise u forwards the packet incrementing the time-to-live counter. If this counter reaches the threshold value, then the packet will be discarded.

The following steps describe the update scenario:

(i) **Preparation Phase:**

The node s_i prepares its new packet with the new and old data along with its ID and counter. Also, s_i add an update counter *token* initialized at 1 for the first updated packet. So, we assume that the following steps happen when *token* is set to 1.

$$packet_{s_i} = (ID_{s_i}, x_{s_i} \oplus x'_{s_i}, c(s_i)). \quad (29)$$

s_i chooses at random a neighbor node u , and sends its $packet_{s_i}$.

(ii) **Encoding Phase:**

The node u checks if the $packet_{s_i}$ is an update or first-time packet. If it is first-time packet it will accept, forward, or discard it as shown in LTCDS-I algorithm 1. If $packet_{s_i}$ is an updated packet, then the node u will check if ID_{s_i} is already included in its accepted list. If yes, then it will update its value y_u as follows.

$$y_u^+ = y_u^- \oplus x_{s_i} \oplus x'_{s_i}. \quad (30)$$

If no, it will add this updated packet into its forward queue with incrementing the counter

$$c(x'_{s_i}) = c(x_{s_i}) + 1. \quad (31)$$

The $packet_{s_i}$ will be discarded if $c(x'_{s_i}) \geq C_1 n \log n$ where C_1 is a system parameter. In this case, we need C_1 to be large enough, so all old data x_{s_i} will be updated to the new data x'_{s_i} .

(iii) **Storage Phase:**

If all nodes are done with updating their values y_i . One can run the decoding phase to retrieve the original and update information.

Now, since we run only one simple random walk for each update, if h is the number of nodes updating their values, then we have the following result.

Lemma 14. *The total number of transmissions needed for the update process is bounded by $\Theta(hn \log n)$.*

5 Performance Evaluation

In this section, we study performance of the proposed LTCDS-I and LTCDS-II algorithms for distributed storage in wireless sensor networks through simulation. The main performance metric we investigate is the *successful decoding probability* versus the *decoding ratio*.

Definition 15. (Decoding Ratio) Decoding ratio η is the ratio between the number of queried nodes h and the number of sources k , i.e.,

$$\eta = \frac{h}{k}. \quad (32)$$

Definition 16. (Successful Decoding Probability) Successful decoding probability P_s is the probability that the k source packets are all recovered from the h querying nodes.

In our simulation, P_s is evaluated as follows. Suppose the network has n nodes and k sources, and we query h nodes. There are $\binom{n}{h}$ ways to choose such h nodes, and we pick one tenth of these choices uniformly at random:

$$M = \frac{1}{10} \binom{n}{h} = \frac{n!}{10 \cdot h!(n-h)!}. \quad (33)$$

Let M_s be the size of the subset these M choices of h query nodes from which the k source packets can be recovered. Then, we evaluate the successful decoding probability as

$$P_s = \frac{M_s}{M}. \quad (34)$$

Figure 3 shows the decoding performance of LTCDS-I algorithm with Ideal Soliton distribution with small number of nodes and sources. The network is deployed in $\mathcal{A} = [5, 5]^2$, and the system parameter C_1 is set as $C_1 = 5$. From the simulation results we can see that when the decoding ratio is above 2, the successful decoding probability is about 99%. Another observation is that when the total number of nodes increases but the ratio between k and n and the decoding ratio η are kept as constants, the successful decoding probability P_s increases when $\eta \geq 1.5$ and decreases when $\eta < 1.5$. This is also confirmed by the results shown in Figure 4. In Figure 4, The network has constant density as $\lambda = \frac{40}{9}$ and the system parameter $C_1 = 3$.

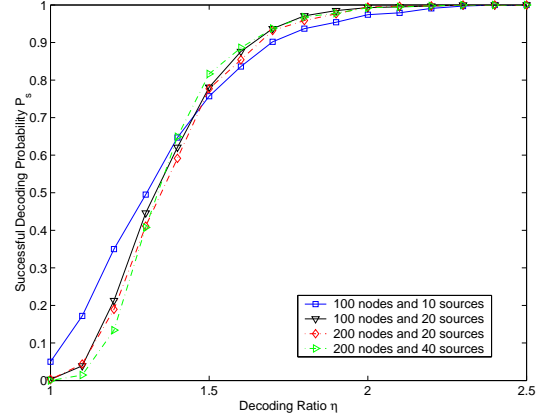


Figure 3. Decoding performance of LTCDS-I algorithm with small number of nodes and sources

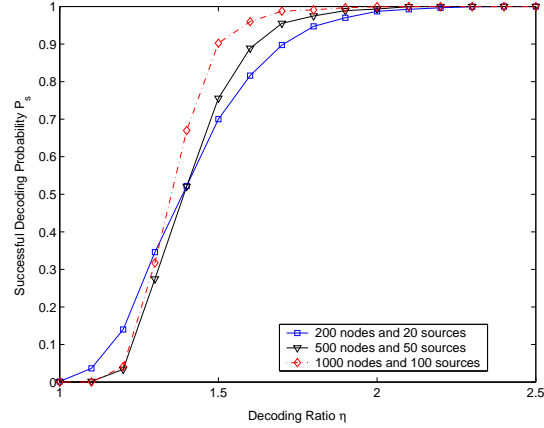


Figure 4. Decoding performance of LTCDS-I algorithm with medium number of nodes and sources

In Figure 5, we fix the decoding ratio η as 1.4 and 1.7, respectively, and fix the ratio between the number of sources and the number of nodes as 10%, i.e., $k/n = 0.1$, and change the number of nodes n from 500 to 5000. From the results, it can be seen that as n grows, the successful decoding probability increases until it reaches some platform which is the successful decoding probability of real LT codes. This confirms that LTCDS-I algorithm has the same asymptotical performance as LT codes.

To investigate how the system parameter C_1 affects the decoding performance of the LTCDS-I algorithm, we fix the decoding ratio η and change C_1 . The simulation results are shown in Figure 6. For the scenario of 1000 nodes and 100 sources, η is set as 1.6, and for the scenario of 500 nodes

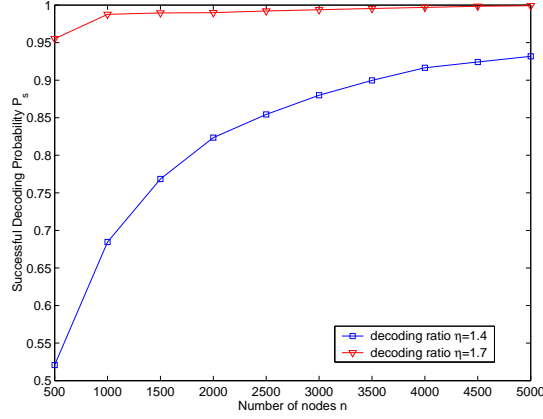


Figure 5. Decoding performance of LTCDS-I algorithm with different number of nodes

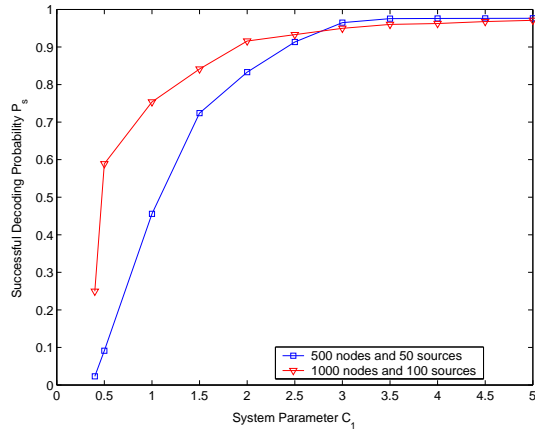


Figure 6. Decoding performance of LTCDS-I algorithm with different system parameter C_1

and 50 sources, η is set as 1.8. The code degree distribution is also the Ideal Soliton distribution, and the network is deployed in $\mathcal{A} = [15, 15]^2$. It can be seen that when $C_1 \geq 3$, P_s keeps almost like a constant, which indicates that after $3n \log n$ steps, almost all source packets visit each node at least once.

Figure 7 compares the decoding performance of LTCDS-II and LTCDS-I with Ideal Soliton distribution with small number of nodes and sources. As in Figure 3, the network is deployed in $\mathcal{A} = [5, 5]^2$, and the system parameter is set as $C_3 = 10$. To guarantee each node obtain accurate estimations of n and k , we set $C_2 = 50$. It can be seen that the decoding performance of the LTCDS-II algorithm is a little bit worse than the LTCDS-I algorithm when decoding ratio η is small, and almost the same when η is large. Figure 8 compares the decoding performance of LTCDS-II and

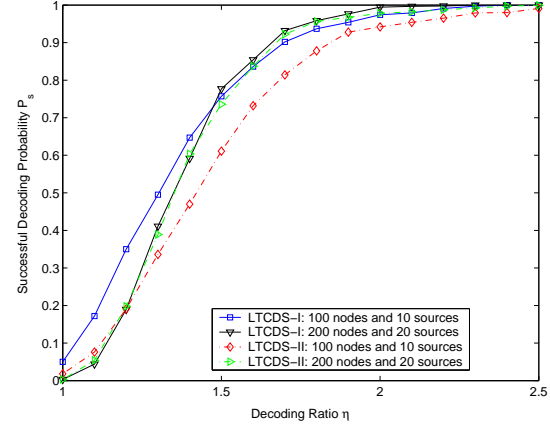


Figure 7. Decoding performance of LTCDS-II algorithm with small number of nodes and sources

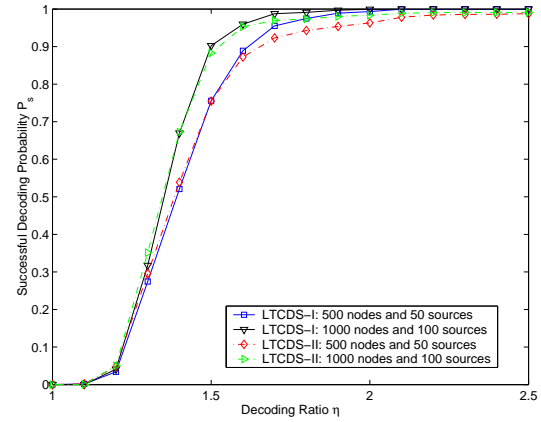


Figure 8. Decoding performance of LTCDS-II algorithm with medium number of nodes and sources

LTCDS-I with Ideal Soliton distribution with medium number of nodes and sources, where the network has constant density as $\lambda = \frac{40}{9}$ and the system parameter $C_3 = 20$. We observe different phenomena. The decoding performance of the LTCDS-II algorithm is a little bit better than the LTCDS-I algorithm when decoding ratio η is small, and almost the same when η is large. That is because for the simulation in Figure 8, we set $C_3 = 20$ which is larger than $C_3 = 10$ set for the simulation in Figure 6. The larger value of C_3 guarantees that each node has the chance to accept each source packet, which results in a more uniformly distribution.

Figure 9–Figure 10 shows the histogram of the estimation results of n and k of each node for three scenarios: Fig-

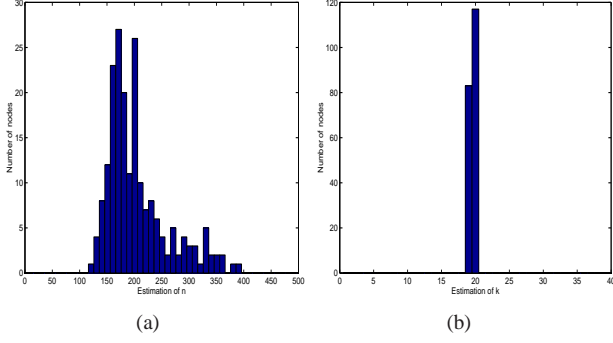


Figure 9. Estimation results in LTCDS-II algorithm with $n = 200$ nodes and $k = 20$ sources: (a) estimations of n ; (b) estimations of k .

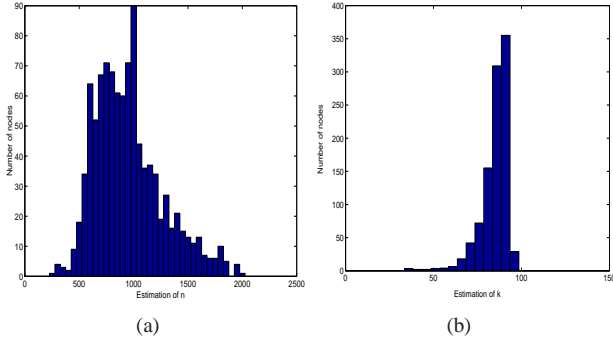


Figure 10. Estimation results in LTCDS-II algorithm with $n = 1000$ nodes and $k = 100$ sources: (a) estimations of n ; (b) estimations of k .

Figure 9 shows the results for 200 nodes and 20 sources; and Figure 10 shows the results for 1000 nodes and 100 sources. In the first two scenarios, we set $C_2 = 50$. From the results we can see that, the estimations of k are more accurate and concentrated than the estimations of n . This is because the estimation of k only depends on the ratio between the expected inter-visit time and the expected inter-packet time, which is independent of the mean degree μ and the node degree $d_n(u)$. On the other hand, the estimation of n actually depends on μ and $d_n(u)$. However, in the LTCDS-II algorithm, each node approximates μ as its own node degree $d_n(u)$, which causes the deviation of the estimations of n .

To investigate how the system parameter C_2 affects the decoding performance of the LTCDS-II algorithm, we fix the decoding ratio η and C_3 , and change C_2 . The simulation results are shown in Figure 11. From the simulation results, we can see that when C_2 is chosen to be small, the

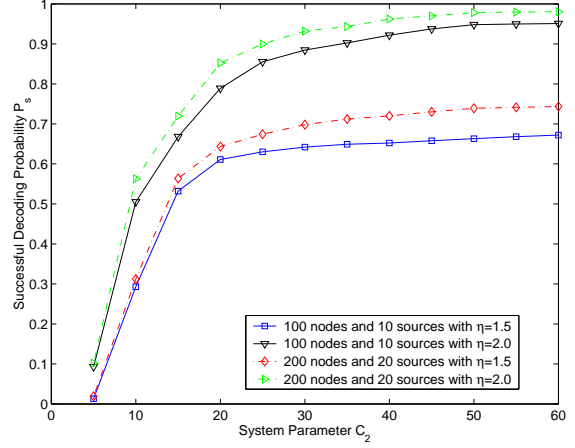


Figure 11. Decoding performance of LTCDS-II algorithm with different system parameter C_2

performance of the LTCDS-II algorithm is very poor. This is due to the inaccurate estimations of k and n of each node. When C_2 is large, for example, when $C_2 \geq 30$, the performance is almost the same.

6 Conclusion

In this paper, we studied a model for large-scale wireless sensor networks, where the network nodes have low CPU power and limited storage. We proposed two new decentralized algorithms that utilize Fountain codes and random walks to distribute information sensed by k sensing source nodes to n storage nodes. These algorithms are simpler, more robust, and less constrained in comparison to previous solutions that require knowledge of network topology, maximum degree of a node, or knowing values of n and k [4, 6, 9, 10, 11]. We computed the computational encoding and decoding complexity of these algorithms and simulated their performance with small and large numbers of k and n nodes. We showed that a node can successfully estimate the number of sources and total number of nodes if it can only compute the *inter-visit time* and *inter-packet time*.

Our future work will include Raptor codes based distributed networked storage algorithms for sensor networks. We also plan to provide theoretical results and proofs for the results shown in this paper, where the limited space is not an issue. Our algorithm for estimating values of n and k is promising, we plan to investigate other network models where this algorithm is beneficial and can be utilized.

Acknowledgments

The authors would like to thank the reviewers for their comments. They would like to express their gratitude to all Bell Labs & Alcatel-Lucent staff members for their hospitality and kindness.

7 Appendix

7.1 Proof of Lemma 10

Proof. For a simple random walk on an undirected graph $G = (V, E)$, the stationary distribution is given by [1, 17, 14]

$$p(u) = \frac{d_n(u)}{2|E|}. \quad (35)$$

On the other hand, for a reversible Markov chain, the expected return time for a state i is given by [1, 17, 14]

$$E[T_{\text{return}}(i)] = \frac{1}{\pi(i)}, \quad (36)$$

where $\pi(i)$ is the stationary distribution of state i .

From (35) and (36), we have for a simple random on a graph, the expected inter-visit time of node u is

$$E[T_{\text{visit}}(u)] = \frac{2|E|}{d_n(u)} = \frac{\mu n}{d_n(u)}, \quad (37)$$

where μ is the mean degree of the graph. \square

7.2 Proof of Lemma 12

Proof. For a given node u and k simple random walks, each simple random walk has expected inter-visit time $\frac{\mu n}{d_n(u)}$. We now view this process from another perspective: we assume there are k nodes $\{v_1, \dots, v_k\}$ uniformly distributed in the network and an agent from node u follows a simple random walk. Then the expected inter-visit time for this agent to visit any particular v_i is the same as $\frac{\mu n}{d_n(u)}$. However, the expected inter-visit time for any two nodes v_i and v_j is $\frac{1}{k} \frac{\mu n}{d_n(u)}$, which gives the expected inter-packet time. \square

References

- [1] D. Aldous and J. Fill. *Reversible Markov Chains and Random Walks on Graphs*. Preprint, available at <http://statwww.berkeley.edu/users/aldous/RWG/book.html>, 2002.
- [2] S. A. Aly, Z. Kong, and E. Soljanin. Fountain codes based distributed storage algorithms. *U.S. patent*, Submitted, October, 2007.
- [3] C. Avin and G. Ercal. On the cover time of random geometric graphs. In *Proc. 32nd International Colloquium of Automata, Languages and Programming, ICALP'05*, Lisboa, Portugal, July, 2005.
- [4] A. G. Dimakis, V. Prabhakaran, and K. Ramchandran. Decentralized erasure codes for distributed networked storage. *IEEE/ACM Transactions on Networking (TON)*, 14(SI):2809 – 2816, June 2006.
- [5] A. G. Dimakis, V. Prabhakaran, and K. Ramchandran. Ubiquitous access to distributed data in large-scale sensor networks through decentralized erasure codes. In *Proc. of 4th IEEE Symposium on Information Processing in Sensor Networks (IPSN '05)*, Los Angeles, CA, USA, April, 2005.
- [6] A. G. Dimakis, V. Prabhakaran, and K. Ramchandran. Distributed fountain codes for networked storage. *Acoustics, Speech and Signal Processing, ICASSP 2006*, may 2006.
- [7] E. N. Gilbert. Random plane networks. *J. Soc. Indust. Appl. Math.*, 9:533–543, 1961.
- [8] A. Kamra, J. Feldman, V. Misra, and D. Rubenstein. Data persistence in sensor networks: Towards optimal encoding for data recovery in partial network failures. In *Workshop on Mathematical performance Modeling and Analysis*, June 2005.
- [9] A. Kamra, V. Misra, J. Feldman, and D. Rubenstein. Growth codes: Maximizing sensor network data persistence. In *Proc. of the 2006 conference on Applications, technologies, architectures, and protocols for computer communications, Sigcomm06*, pages 255 – 266, Pisa, Italy, 2006.
- [10] Y. Lin, B. Li, , and B. Liang. Differentiated data persistence with priority random linear code. In *Proc. of 27th International Conference on Distributed Computing Systems (ICDCS'07)*, Toronto, Canada, June, 2007.
- [11] Y. Lin, B. Liang, and B. Li. Data persistence in large-scale sensor networks with decentralized fountain codes. In *Proc. of the 26th IEEE INFOCOM07*, Anchorage, Alaska, May 6-12, 2007.
- [12] M. Luby. LT codes. In *Proc. 43rd Symposium on Foundations of Computer Science (FOCS 2002)*, 16–19 November 2002, Vancouver, BC, Canada, 2002.
- [13] D. S. Lun, N. Ranakar, R. Koetter, M. Medard, E. Ahmed, and H. Lee. Achieving minimum-cost multicast: A decentralized approach based on network coding. In *In Proc. the 24th IEEE INFOCOM*, volume 3, pages 1607– 1617, March 2005.
- [14] R. Motwani and P. Raghavan. *Randomized Algorithms*. Cambridge University Press, 1995.
- [15] M. Penrose. *Random Geometric Graphs*. Oxford University Press, New York, 2003.
- [16] M. Pitkanen, R. Moussa, M. Swamy, and T. Niemi. Erasure codes for increasing the availability of grid data storage. In *Proc. of the Advanced International Conference on Telecommunications and International Conference on Internet and Web Applications and Services (AICT/ICIW)*, 2006.
- [17] S. Ross. *Stochastic Processes*. Wiley, New York, second edition, 1995.
- [18] I. Stojmenovic. *Handbook of sensor networks, algorithms and architectures*. *Wiley series on parallel and distributed computing*, 2005.